# RUC_AIM3 at TRECVID 2020: Ad-hoc Video Search

**Yida Zhao**, Yuqing Song, Shizhe Chen, Qin Jin*

**AI·M³ , Renmin University of China**

➢ We learn, think, and express through multiple modalities
➢ AI system needs to have the ability to understand the multimodal world
➢ We focus on understanding from Multi-Level, Multi-Aspect, and Multi-Modal (M³)

MULTI ASPECT MULTI MODAL MULTI LEVEL

$AI \cdot M^3$

中国人民大学多媒体计算实验室
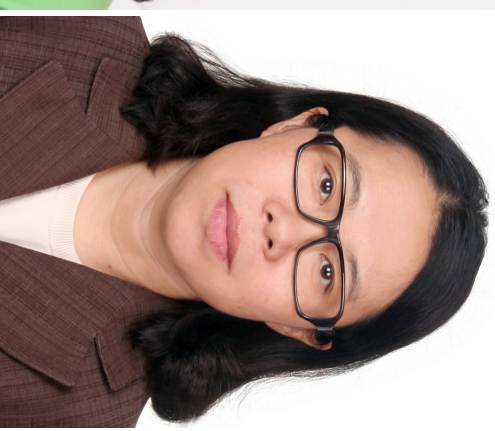
## RUC_AIM3 Team

Yida Zhao

Yuqing Song

Shizhe Chen

Qin Jin

# Outline

- Task Introduction
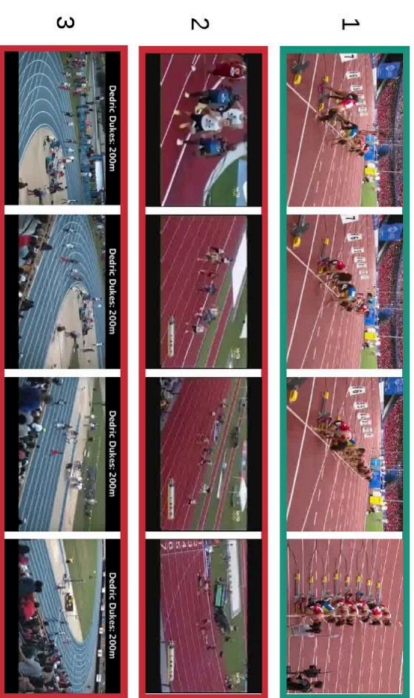
- The Proposed System

- Conclusions

# One

# Task Introduction

# Task Introduction

- ## Ad-hoc Video Search

Given a text query, retrieve the most relevant top 1000 video clips from the V3C vimeo collection [1], which contains about one million video clips

Text: several woman are setting up in the blocks preparing to start a track race.

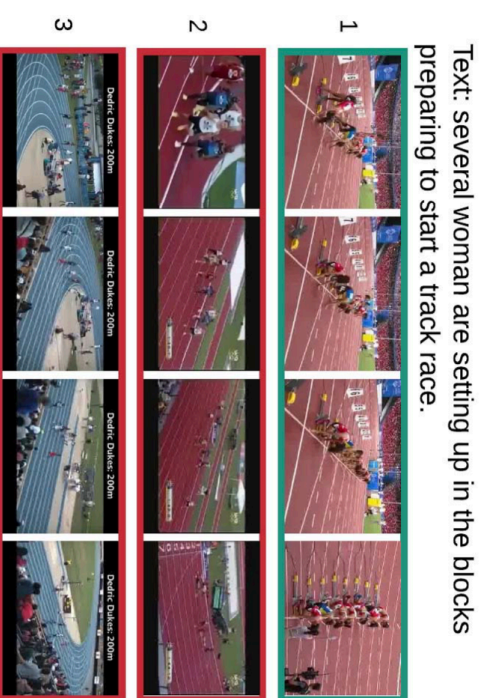[1] V3c–a research video collection, ICMM, 2019

# Task Introduction

- Ad-hoc Video Search

Given a text query, retrieve the most relevant top 1000 video clips from the V3C vimeo collection [1], which contains about one million video clips
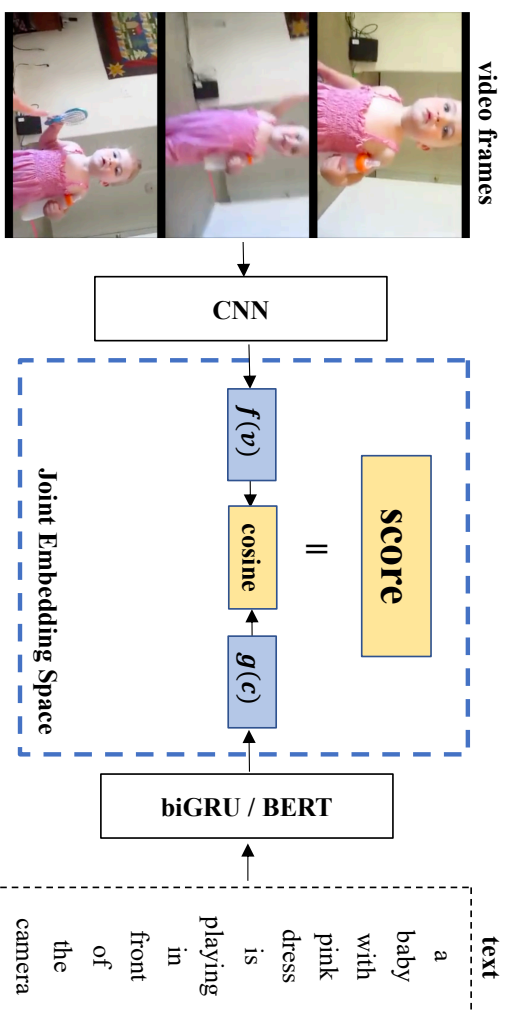
- Challenge

The semantic matching between videos and texts

[1] V3c–a research video collection, ICMM, 2019
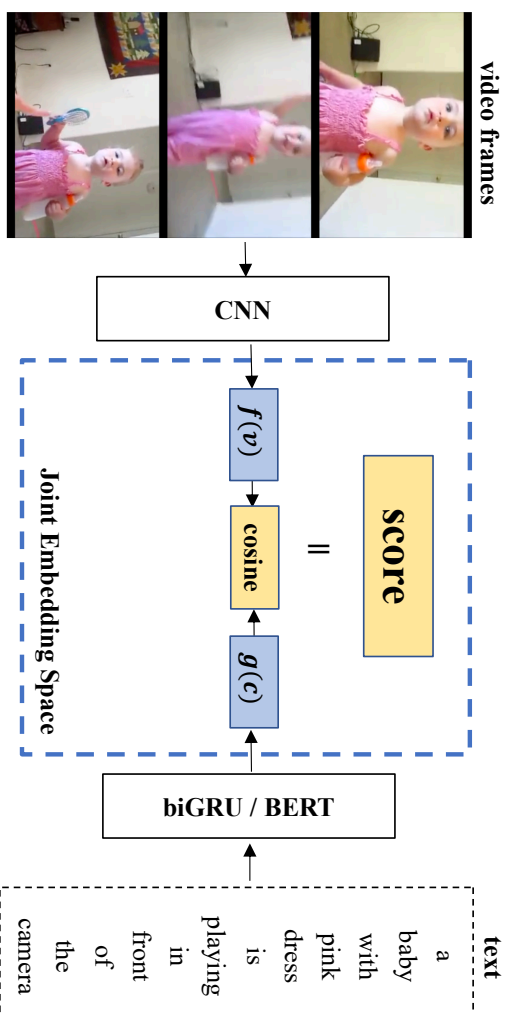
Text: several woman are setting up in the blocks preparing to start a track race.

# Task Introduction

- Video-Text Cross-modal Retrieval
  - Dominant approach: learning joint embedding space and global visual-semantic matching

# Task Introduction

- Video-Text Cross-modal Retrieval

- Dominant approach: learning joint embedding space and global visual-semantic matching

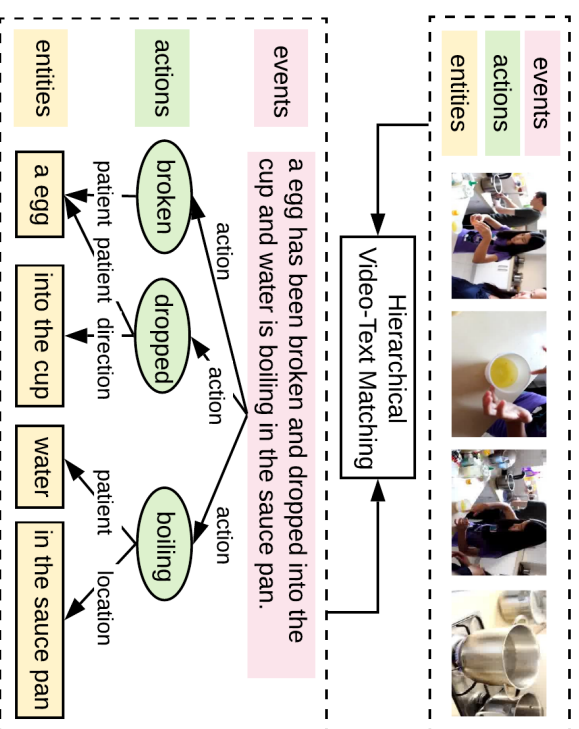video frames

CNN

$f(v)$

cosine

$g(c)$

**score**

=

Joint Embedding Space

biGRU / BERT

text

a
baby
with
pink
dress
is
playing
in
front
of
the
camera

☹ One vector is hard to encode fine-grained details
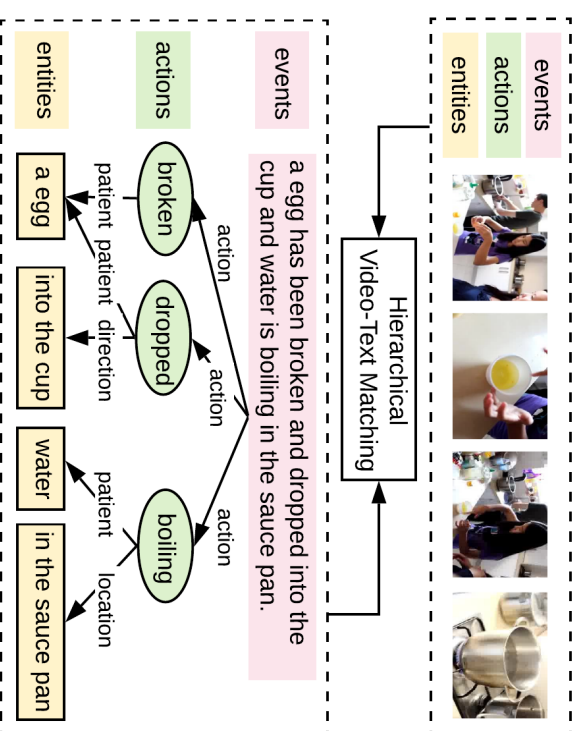
*Two*

# Our System

## Two-branch Model

- Two-branch Matching Model
  - Global Matching
    - VSE++ [2]
    - Dual Encoding [3]
  - Fine-grained Matching
    - HGR [4]

[2] Vse++: Improving visual-semantic embeddings with hard negatives, BMVC, 2018
[3] Dual encoding for zero-example video retrieval, CVPR, 2019
[4] Fine-grained video-text retrieval with hierarchical graph reasoning, CVPR, 2020

- Two-branch Matching Model

  - **Global Matching**

    - **VSE++ [2]**

    - **Dual Encoding [3]**

  - Fine-grained Matching

    - HGR [4]

[2] Vse++: Improving visual-semantic embeddings with hard negatives, BMVC, 2018

[3] Dual encoding for zero-example video retrieval, CVPR, 2019

[4] Fine-grained video-text retrieval with hierarchical graph reasoning, CVPR, 2020

- Two-branch Matching Model
  - Global Matching
    - VSE++ [2]
    - Dual Encoding [3]
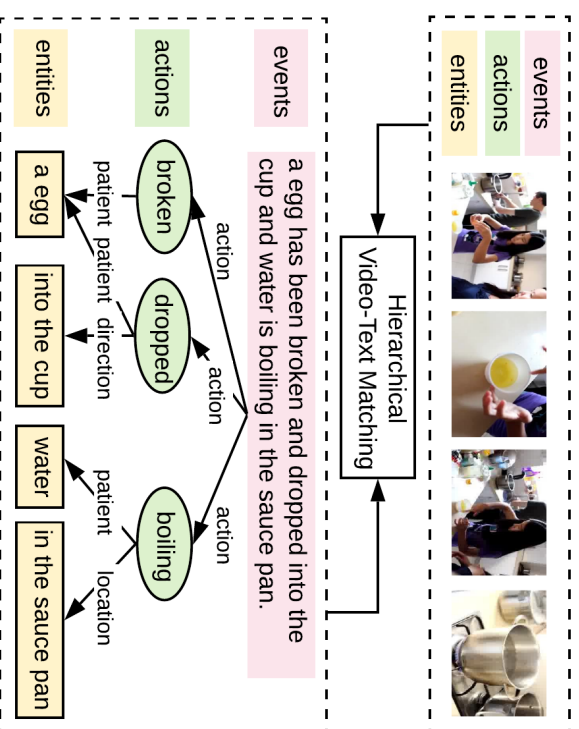  - **Fine-grained Matching**
    - **HGR [4]**



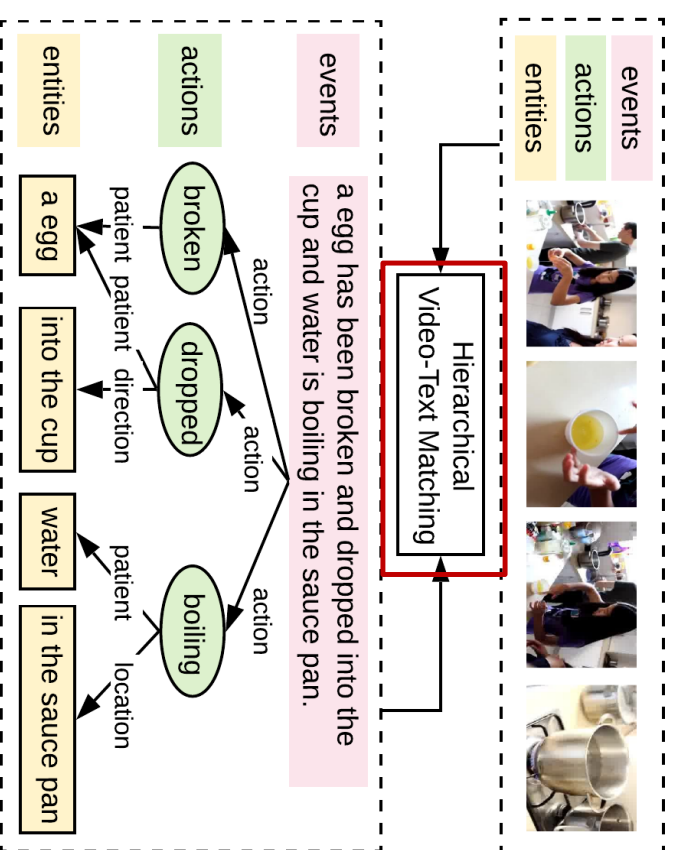[2] Vse++: Improving visual-semantic embeddings with hard negatives, BMVC, 2018
[3] Dual encoding for zero-example video retrieval, CVPR, 2019
[4] Fine-grained video-text retrieval with hierarchical graph reasoning, CVPR, 2020

- Multi-level Video-Text Matching
  - Event
  - Actions
  - Entities

Global

Local

events
actions
entities

Hierarchical Video-Text Matching

a egg has been broken and dropped into the cup and water is boiling in the sauce pan.

events

actions

entities

broken
dropped
boiling

action
action
action

patient
patient
direction
patient
location

a egg
into the cup
water
in the sauce pan

- Multi-level Video-Text Matching
  - Event
  - Actions
  - Entities

  Global ↓ Local

- Hierarchical Textual Encoding
  - Decompose sentence into semantic role graph
  - Capture relationships via graph reasoning

- Textual Graph Construction
  - Event node: the whole text query
  - Action node: verbs in the text
  - Entity node: noun phrases in the text

# Hierarchical Graph Reasoning (HGR)

- Textual Graph Construction
  - Event node: the whole text query
  - Action node: verbs in the text
  - Entity node: noun phrases in the text
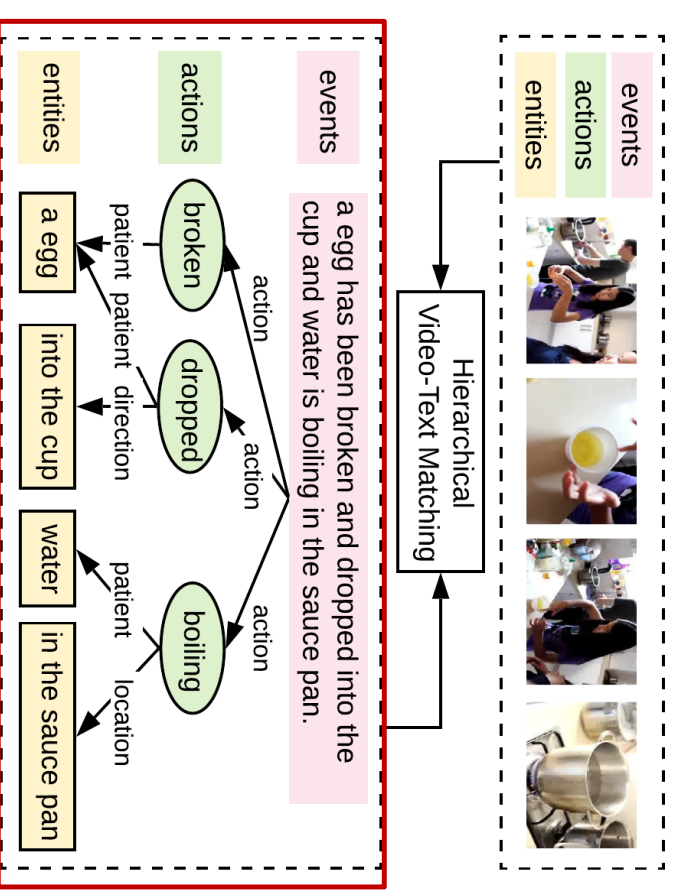- Attentive Graph Reasoning —— Relational GCN

# Hierarchical Graph Graph Reasoning (HGR)

- Multi-level Video-Text Matching
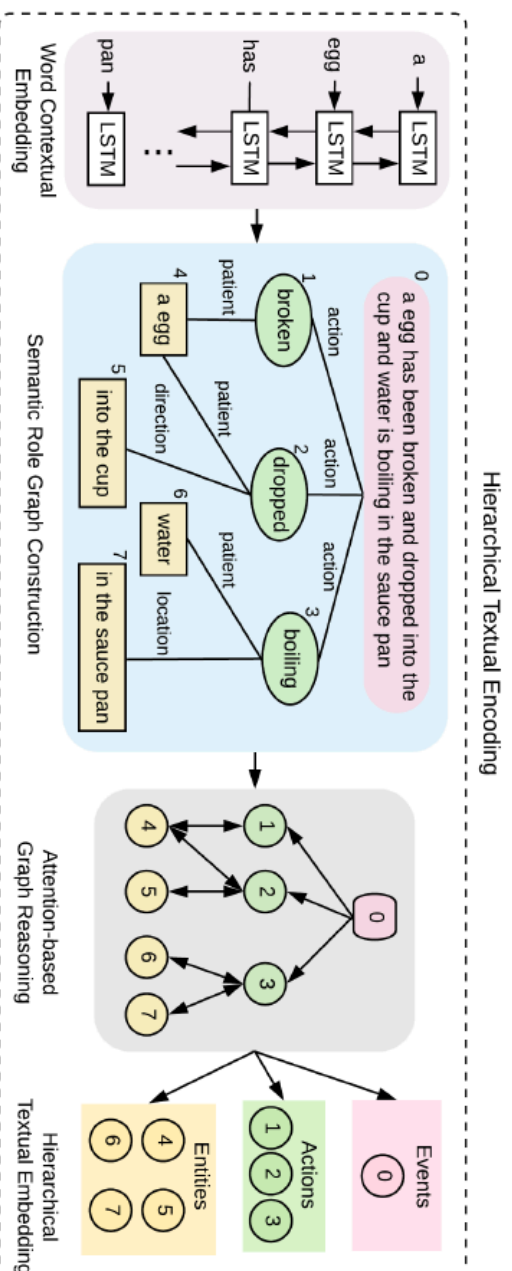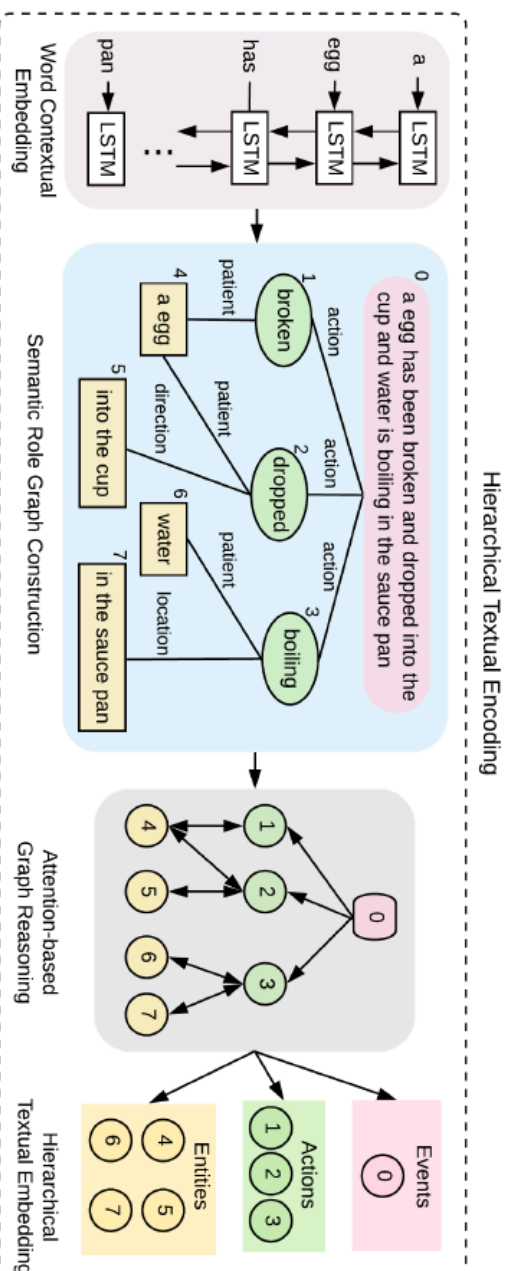  - Event
  - Actions
  - Entities

  **Global** | **Local**

- Hierarchical Textual Encoding
  - Decompose sentence into semantic role graph
  - Capture relationships via graph reasoning

- Hierarchical Video Encoding
  - Guided by different levels of text to learn diverse video representations

- Video Encoding

$$V = \{f_1, ..., f_M\}$$

$$v_{x,i} = W_x^v f_i, \quad x \in \{e, a, o\}$$



Hierarchical Video Encoding

Events

Actions

Entities

Hierarchical
Video Embedding

- Global Matching at the event level

$$s_e = cos(v_e, c_e)$$

- Local Attentive Matching

$$s_{ij}^x = cos(v_{x,j}, c_{x,i})$$

$$s_{x,i} = \sum_j \varphi_{ij}^x s_{ij}^x$$

$$\varphi_{ij}^x = \text{softmax}(\lambda[s_{ij}^x]_+ / \sqrt{\sum_j [s_{ij}^x]_+^2})$$

$$s_x = \sum_i s_{x,i} \quad x \in \{a, o\}$$

- $s(v, c) = (s_e + s_a + s_o)/3$



Video-Text Matching

Global Match

Local
Attentive Match
(Action & Entity)

# Video-Text Matching Results

- **Video datasets** : TGIF, MSRVTT, VATEX

  **Image dataset** : MSCOCO (only for global matching models)

- **Video features** : ResNeXt-101, irCSN-152

  **Image features** : ResNeXt-101

# Video-Text Matching Results

- Four runs for the final submission:
  - Run4: The global matching branch trained on video datasets

Table1. Results on TRECVID 2019 and 2020 AVS Main Task.

| Submissions | 2019 | 2020 |
|---|---|---|
| Winner in 2019 | 0.163 | - |
| Run4 | 0.177 | 0.354 |

# Video-Text Matching Results

- Four runs for the final submission:
  - Run4: The global matching branch trained on video datasets
  - Run3: Run4 + global matching branch trained on image datasets

Table1. Results on TRECVID 2019 and 2020 AVS Main Task.

| Submissions | 2019 | 2020 |
|---|---|---|
| Winner in 2019 | 0.163 | - |
| Run4 | 0.177 | 0.354 |
| Run3 | 0.193 | 0.350 |

# Video-Text Matching Results

- Four runs for the final submission:
  - Run4: The global matching branch trained on video datasets
  - Run3: Run4 + global matching branch trained on image datasets
  - Run2: Run3 + fine-grained matching branch (HGR)

**Table1. Results on TRECVID 2019 and 2020 AVS Main Task.**

| Submissions | 2019 | 2020 |
|---|---|---|
| Winner in 2019 | 0.163 | - |
| Run4 | 0.177 | 0.354 |
| Run3 | 0.193 | 0.350 |
| Run2 | 0.195 | 0.357 |

# Video-Text Matching Results

- Four runs for the final submission:
  - Run4: The global matching branch trained on video datasets
  - Run3: Run4 + global matching branch trained on image datasets
  - Run2: Run3 + fine-grained matching branch (HGR)
  - Run1: Run2 + BERT as text encoder

**Table1. Results on TRECVID 2019 and 2020 AVS Main Task.**

| Submissions | 2019 | 2020 |
|---|---|---|
| Winner in 2019 | 0.163 | - |
| Run4 | 0.177 | 0.354 |
| Run3 | 0.193 | 0.350 |
| Run2 | 0.195 | 0.357 |
| **Run1** | **0.196** | **0.359** |

# Video-Text Matching Results

- Four runs for the final submission:

  - Run4: The global matching branch trained on video datasets
  - Run3: Run4 + global matching branch trained on image datasets
  - Run2: Run3 + fine-grained matching branch (HGR)
  - Run1: Run2 + BERT as text encoder

**Table1. Results on TRECVID 2019 and 2020 AVS Main Task.**

| Submissions | 2019 | 2020 |
|---|---|---|
| Winner in 2019 | 0.163 | - |
| Run4 | 0.177 | 0.354 |
| Run3 | 0.193 | 0.350 |
| Run2 | 0.195 | 0.357 |
| **Run1** | **0.196** | 0.359 |
| **Run5**\* | 0.181 | **0.361** |

# Video-Text Matching Results

- Four runs for the final submission:
  - Run4: The global matching branch trained on video datasets
  - Run3: Run4 + global matching branch trained on image datasets
  - Run2: Run3 + fine-grained matching branch (HGR)
  - Run1: Run2 + BERT as text encoder

**Table2. Results on TRECVID AVS Progress Subtask.**

| Submissions | Results |
|---|---|
| Winner in 2019 | 0.177 |
| **Run4** | **0.235** |
| Run3 | 0.208 |
| Run2 | 0.220 |
| Run1 | 0.223 |

# Take Home Message

- We propose a two-branch model by combing the global matching and fine-grained matching for the AVS task

# Take Home Message

- We propose a two-branch model by combing the global matching and fine-grained matching for the AVS task

- Training on additional image captioning dataset can improve the retrieval performance on 2019 AVS task, but not on 2020 AVS task

# Take Home Message

- We propose a two-branch model by combing the global matching and fine-grained matching for the AVS task

- Training on additional image captioning dataset can improve the retrieval performance on 2019 AVS task, but not on 2020 AVS task

- Our models rank the 1st place on the TRECVID 2020 AVS Main Task.

# **THANKS !**

If you have any questions , please feel free to contact with us:

zyiday@ruc.edu.cn, syuqing@ruc.edu.cn, cszhe1@ruc.edu.cn, qjin@ruc.edu.cn

http://jin-qin.com/AIM3-Lab.html